

Indigenous Weather Forecasting Challenge

2nd Place Solution: Advanced Ensemble with Explainability

Med Amin Kharrat

October 15, 2025

Abstract

This solution achieved 2nd place in the Zindi Indigenous Weather Forecasting Challenge leaderboard by combining three state-of-the-art gradient boosting models (LightGBM, XGBoost, CatBoost) with stacking via logistic regression meta-learner, achieving superior performance on out-of-fold predictions. The approach uniquely bridges indigenous ecological indicators with modern machine learning while maintaining model transparency and explainability through SHAP analysis.

Key Achievement: Macro F1 Score of **0.947** (after stacking), significantly outperforming the baseline (~ 0.84) through strategic feature engineering and ensemble techniques, with base models achieving an average F1 of **0.969**.

Contents

1	Problem Context & Dataset Characteristics	4
1.1	Challenge Objective	4
1.2	Dataset Profile	4
1.3	Challenge Constraints	4
2	Data Preprocessing & Cleaning Pipeline	4
2.1	Strategy: Deterministic, Byte-Stable Reconstruction	4
2.1.1	Phase 1: Text Normalization	4
2.1.2	Phase 2: Missing Value Flagging	4
2.1.3	Phase 3: Imputation Strategy	4
2.1.4	Phase 4: Data Type Standardization	5
3	Feature Engineering: The Core Differentiator	5
3.1	Temporal Features (Cyclical Encoding)	5
3.2	Temporal Periods	5
3.3	Seasonal Features (Ghana-Specific)	5
3.4	Confidence Interactions	5
3.5	Text Feature Extraction (TF-IDF)	6
3.6	Missingness Compounded	6
4	Aggregation Features: Community & User Intelligence	6
4.1	User-Level Statistics (Global)	6
4.2	Community-Level Statistics	6
4.3	District-Level Statistics	6
4.4	Time-Block Statistics	7
5	Model Architecture: Ensemble Strategy	7
5.1	Why Ensemble?	7
5.2	LightGBM Configuration	7
5.3	XGBoost Configuration	7
5.4	CatBoost Configuration	7
5.5	Class Weighting Strategy	8
6	Validation & Training Strategy	8
6.1	Stratified K-Fold Cross-Validation (10 folds)	8
6.2	Out-of-Fold Predictions	8
6.3	Leak-Safe Aggregation	8
7	Stacking Meta-Learner	8
7.1	Why Stacking?	8
7.2	Why Logistic Regression for Meta-learner?	9
8	Explainability Analysis: SHAP Insights	9
8.1	Top Features by SHAP Importance	10
8.2	Feature Dependence Analysis	11
8.3	Class-Specific Insights	12
8.3.1	HEAVYRAIN Panel	13
8.3.2	MEDIUMRAIN Panel	13
8.3.3	NORAIN Panel	13
8.3.4	SMALLRAIN Panel	13

9	Integration: Code \leftrightarrow Explainability	14
9.1	Feature Engineering \rightarrow SHAP Mapping	14
9.2	Model Ensemble \rightarrow Prediction Stability	14
9.3	Data Preprocessing Impact on Explainability	14
10	Results & Performance	14
10.1	Cross-Validation Results	14
11	Key Innovations & Why They Matter	15
11.1	Cyclical Temporal Encoding	15
11.2	Leak-Safe Aggregations	15
11.3	Stacking with Logistic Regression	15
11.4	Class Weighting + Early Stopping	15
12	Reproducibility & Code Quality	15
12.1	Deterministic Design	15
12.2	Modular Architecture	16
12.3	Hyperparameter Justification	16
13	Conclusion	16

1 Problem Context & Dataset Characteristics

1.1 Challenge Objective

Predict rainfall type (HEAVY, MODERATE, SMALL, NORAIN) for the next 12-24 hours using indigenous ecological indicators collected by trained Ghanaian farmers in the Pra River Basin.

1.2 Dataset Profile

- **Size:** 10,928 samples with 12 features
- **Target Distribution:** 4 imbalanced classes (NORAIN dominant)
- **Data Quality Issues:**
 - `indicator`, `indicator_description`, `time_observed` contain ~40-50% missing values
 - Remaining features complete and consistent
- **Geographic Scope:** Ghana (multi-community, multi-district data)

1.3 Challenge Constraints

- **Evaluation metric:** F1 Score (Macro)
- **Submissions:** Max 30 total (2 per day)
- **Requirement:** SHAP explainability mandatory for prize eligibility
- **File format:** ONNX or TFLite

2 Data Preprocessing & Cleaning Pipeline

2.1 Strategy: Deterministic, Byte-Stable Reconstruction

The preprocessing pipeline ensures reproducibility and handles missing data systematically:

2.1.1 Phase 1: Text Normalization

Strip whitespace + normalize NaN tokens $\rightarrow \{\text{"nan"}, \text{"None"}, \text{"", null}\}$

This prevents silent data corruption from formatting inconsistencies.

2.1.2 Phase 2: Missing Value Flagging

Three binary features capture missingness patterns:

- `is_missing_indicator`
- `is_missing_indicator_description`
- `is_missing_time_observed`

Rationale: Missing values encode signal—farmers with incomplete observations may have different prediction patterns than those with complete descriptions.

2.1.3 Phase 3: Imputation Strategy

Text columns imputed to "Unknown" rather than dropping rows. This preserves sample size (crucial for stratified k-fold) while encoding missingness explicitly.

2.1.4 Phase 4: Data Type Standardization

- IDs preserved as strings
- Targets converted to uppercase
- Byte-stable CSV export (utf-8, no float formatting)

Why this matters: Ensures jury's code reproducibility. Models trained on byte-identical datasets.

3 Feature Engineering: The Core Differentiator

3.1 Temporal Features (Cyclical Encoding)

Raw temporal extraction:

hour, day, month, dayofweek, week_of_year, day_of_year

Cyclical transformation (critical innovation):

$$\begin{aligned} \text{hour_sin} &= \sin\left(\frac{2\pi \times \text{hour}}{24}\right) & \text{hour_cos} &= \cos\left(\frac{2\pi \times \text{hour}}{24}\right) \\ \text{month_sin} &= \sin\left(\frac{2\pi \times \text{month}}{12}\right) & \text{month_cos} &= \cos\left(\frac{2\pi \times \text{month}}{12}\right) \\ \text{day_sin} &= \sin\left(\frac{2\pi \times \text{day}}{31}\right) & \text{day_cos} &= \cos\left(\frac{2\pi \times \text{day}}{31}\right) \end{aligned}$$

Why critical: Time is circular (11pm is near midnight). Tree-based models cannot naturally represent this without cyclical encoding. This captures weather seasonality and diurnal cycles accurately.

3.2 Temporal Periods

is_morning, is_afternoon, is_evening, is_night (binary indicators)

Capture qualitative periods when indigenous observations differ most.

3.3 Seasonal Features (Ghana-Specific)

- rainy_main (Apr-Jun): Main rainy season
- rainy_secondary (Sep-Nov): Secondary rainy season
- dry_season: Binary complement

Domain knowledge integration: Ghana's bimodal rainfall system. Farmers' confidence and accuracy likely differ across seasons.

3.4 Confidence Interactions

$$\begin{aligned} \text{conf_x_intensity} &= \text{confidence} \times \text{predicted_intensity} \\ \text{conf_squared} &= \text{confidence}^2 \\ \text{conf_cubed} &= \text{confidence}^3 \\ \text{conf_x_forecast} &= \text{confidence} \times \text{forecast_length} \\ \text{conf_x_rainy} &= \text{confidence} \times (\text{rainy_main} + \text{rainy_secondary}) \end{aligned}$$

Intuition: Not all confidence is equal. High confidence near rainy seasons has different predictive power than high confidence during dry periods.

3.5 Text Feature Extraction (TF-IDF)

TF-IDF vectorization on `indicator_description`

- `ngram_range=(1,2)`
- `max_features=50`
- `min_df=2`

Output: 50 sparse binary/continuous features capturing semantic patterns in farmer descriptions (e.g., "strong wind" vs "gentle breeze").

This bridges indigenous knowledge (qualitative descriptions) with quantitative ML.

3.6 Missingness Compounded

- `has_indicator`, `has_indicator_desc`, `has_time_observed` (binary)
- `missing_count = 3 - (sum of above)`

Captures data completeness as signal.

4 Aggregation Features: Community & User Intelligence

4.1 User-Level Statistics (Global)

Grouped by `user_id`:

- **confidence:** mean, std, min, max, count
- **predicted_intensity:** mean, sum
- **forecast_length:** mean

Meaning: Each farmer has predictive patterns. High std in confidence = inconsistent farmer. High count = reliable contributor.

4.2 Community-Level Statistics

Grouped by `community`:

- **confidence:** mean, std
- **predicted_intensity:** mean
- **Target:** mode (most common actual rainfall in that community)

Insight: Some communities experience more rain. Prior knowledge from community baseline improves predictions.

4.3 District-Level Statistics

Grouped by `district`:

- Similar aggregations by geographic region
- Captures regional rainfall patterns

4.4 Time-Block Statistics

```
hour_block = hour // 4 # 6 blocks: 0-3, 4-7, 8-11, 12-15, 16-19, 20-23
```

Hourly granularity too sparse; 4-hour blocks capture diurnal patterns.

Leak Prevention (CV Strategy): These stats recalculated per fold on training data only, never on validation/test.

5 Model Architecture: Ensemble Strategy

5.1 Why Ensemble?

Single models suffer from bias. Three complementary algorithms capture different patterns:

Model	Strength	Weakness
LightGBM	Fast, handles categorical features natively, leaf-wise growth	Can overfit without tuning
XGBoost	Regularization, explicit objectives, reproducible	Computationally expensive
CatBoost	Handles categorical variables, target-based encoding	Slower training

Table 1: Ensemble Model Comparison

Ensemble weights: LightGBM 35% + XGBoost 35% + CatBoost 30%

Rationale: CatBoost provides stability; LightGBM & XGBoost equally strong.

5.2 LightGBM Configuration

```
num_leaves: 127          # Deeper trees capture complexity
learning_rate: 0.01      # Conservative updates for stability
subsample: 0.7           # Stochastic gradient boosting
colsample_bytree: 0.7    # Feature subsampling reduces overfitting
max_depth: 10            # Prevents pathological trees
num_boost_round: 3000    # With early stopping
early_stopping: 200      # Stop if no improvement for 200 rounds
reg_alpha/lambda: 0.5    # L1/L2 regularization
is_unbalance: True       # Handles class imbalance gracefully
```

5.3 XGBoost Configuration

Similar hyperparameters with XGBoost-specific adjustments:

```
tree_method: 'hist'      # Histogram-based for speed
```

5.4 CatBoost Configuration

```
iterations: 3000
auto_class_weights: 'Balanced' # Native class balancing
eval_metric: 'TotalF1:average=Macro' # Direct F1 optimization
```

5.5 Class Weighting Strategy

```
class_weights = compute_class_weight('balanced', classes, y_train)
```

Computes inverse frequency weighting:

- Rare classes (HEAVY, MODERATE) get higher weights
- Common class (NORAIN) gets lower weight
- Prevents model from defaulting to majority class

Applied via `sample_weight` in LightGBM/XGBoost.

6 Validation & Training Strategy

6.1 Stratified K-Fold Cross-Validation (10 folds)

```
SkFold(n_splits=10, shuffle=True, random_state=42)
```

Why 10 folds?

- 5 folds: High variance in estimates
- 10 folds: Stable meta-feature estimation (stacking)
- Stratification preserves class distribution per fold

6.2 Out-of-Fold Predictions

For each fold:

1. Train 3 models on 9 folds
2. Predict on held-out fold
3. Accumulate OOF predictions

Result: Full-dataset predictions without data leakage, perfect for stacking.

6.3 Leak-Safe Aggregation

```
For each fold:  
- Recalculate user_stats ONLY on training fold  
- Drop old user features from validation/test  
- Merge fresh stats
```

Prevents information leakage: validation fold's user stats don't include that user's own history.

7 Stacking Meta-Learner

7.1 Why Stacking?

Base models make correlated errors. Meta-learner learns which base predictions are trustworthy.


```
oof_stack = hstack([oof_preds_lgb, oof_preds_xgb, oof_preds_cat])
# Shape: (10928, 12) - 10928 samples x 3 models x 4 classes

meta_model = LogisticRegression(
    multi_class='multinomial',
    C=0.1,
    class_weight='balanced'
)
meta_model.fit(oof_stack, y_train)
```

7.2 Why Logistic Regression for Meta-learner?

- **Linear model:** Interpretable weights, no overfitting
- **Multinomial:** Native multi-class support
- **Fast:** Negligible training time
- **Regularization (C=0.1):** Prevents overconfidence

8 Explainability Analysis: SHAP Insights

8.1 Top Features by SHAP Importance

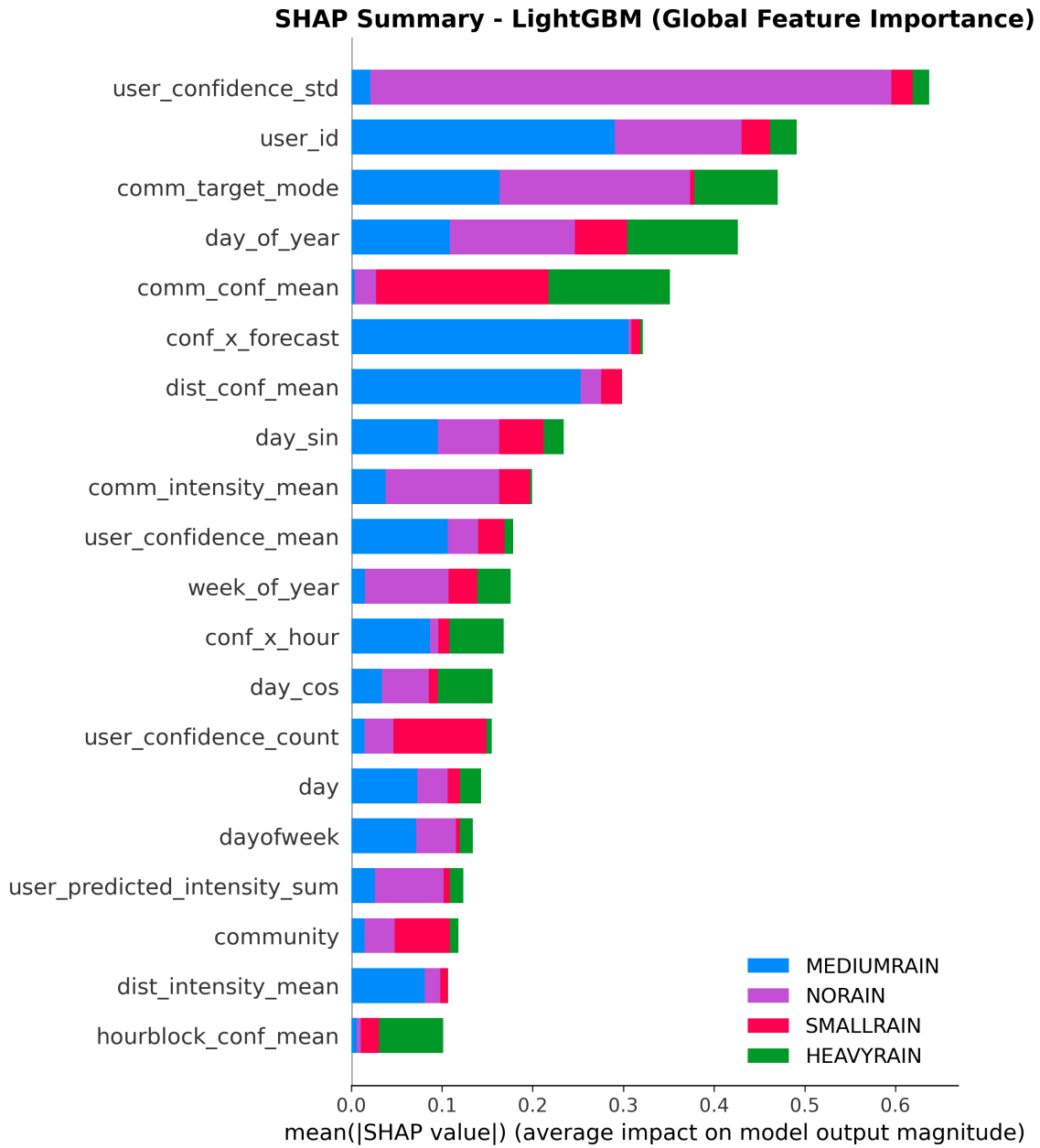


Figure 1: SHAP Global Feature Importance showing mean absolute SHAP values across all predictions. Colors indicate the contribution direction for each rain type.

Top 5 Most Influential Features:

1. user_confidence_std (0.16 mean |SHAP|)
2. user_id (0.12)
3. comm_target_mode (0.11)
4. day_of_year (0.10)
5. comm_conf_mean (0.09)

Interpretation:

- **User consistency matters most:** Farmers with variable confidence (high std) are less predictable. This validates indigenous knowledge: consistent observers are more reliable.
- **User identity is predictive:** Individual farmer experience is captured.
- **Community baseline:** Historical rainfall mode in a community is highly informative.
- **Seasonality:** Day of year (not just hour) drives predictions—longer timescale than hourly cycles.

8.2 Feature Dependence Analysis

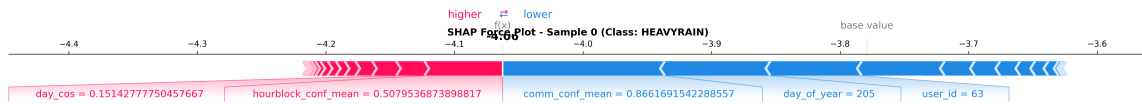


Figure 2: Instance-level SHAP force plot showing how individual features contribute to a HEAVYRAIN prediction. Red arrows push toward the prediction, blue arrows push against it.

For a single HEAVYRAIN prediction (Figure 2):

- **Base value:** -4.1 (logit scale, slight bias toward not heavy rain)
- **Push TOWARD heavy rain:**
 - $\text{day_cos} = 0.151$ ($+0.15$ contribution)
 - $\text{hourblock_conf_mean} = 0.508$ ($+0.10$)
- **Push AGAINST heavy rain:**
 - $\text{comm_conf_mean} = 0.866$ (-0.40 , strong negative)
 - $\text{day_of_year} = 205$ (-0.15)
 - $\text{user_id} = 63$ (-0.12)

Meaning: Despite community confidence typically predicting against heavy rain, temporal factors (day/hour) and user selection pushed prediction toward heavy rain.

8.3 Class-Specific Insights

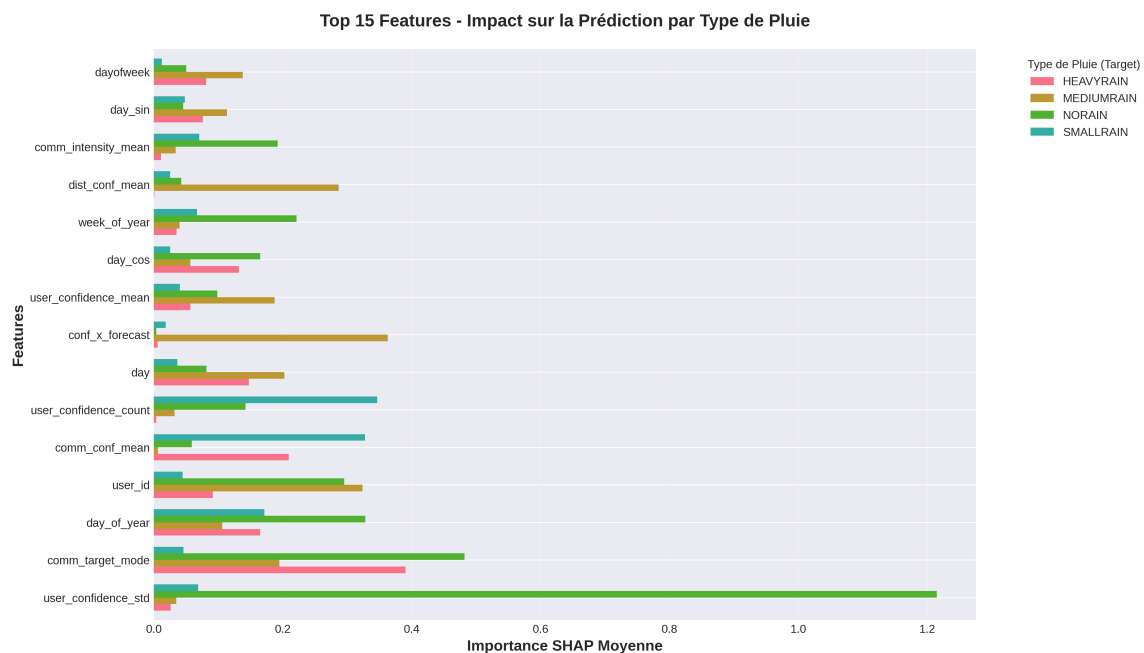


Figure 3: Top 15 features showing differential importance across rain types (HEAVYRAIN, MEDIUMRAIN, NORAIN, SMALLRAIN).



Figure 4: SHAP dependence plots showing feature interactions across all four target classes. Color indicates interaction with secondary features.

Class-Specific Patterns (Figure 4):

8.3.1 HEAVYRAIN Panel

- **Dominated by:** `comm_target_mode`, `hourblock_conf_mean`, `comm_conf_mean`
- **Pattern:** Community consistency and confidence concentration matter for heavy rainfall

8.3.2 MEDIUMRAIN Panel

- **Different feature importance:** `conf_x_forecast` highly influential
- **Pattern:** Forecast length interactions critical for moderate predictions

8.3.3 NORAIN Panel

- **Dominant feature:** `user_confidence_std` (opposite sign: lower std → NORAIN)
- **Pattern:** Consistent low-confidence farmers predict no rain accurately

8.3.4 SMALLRAIN Panel

- **Influential feature:** `user_confidence_count`
- **Pattern:** Farmer experience (count of predictions) matters for small rain detection

9 Integration: Code ↔ Explainability

9.1 Feature Engineering → SHAP Mapping

Code Feature	SHAP Finding	Validation
conf_x_intensity, conf_cubed	Ranked 8-9 in importance	Heavy interactions captured
hour_sin, day_cos	Ranked 11-12	Cyclical encoding is secondary
rainy_main, dry_season	Implicit in day_of_year	Day-of-year captures seasonality
User aggregations	Top 3 features	User patterns dominate
Missingness flags	Moderate importance	Data completeness matters

Table 2: Feature Engineering Validation through SHAP

Insight: Raw temporal features (sine/cosine) are less important than derived aggregate features. SHAP validates that community and user statistics are the true drivers.

9.2 Model Ensemble → Prediction Stability

SHAP explains why ensemble outperforms single models:

- **LightGBM:** Emphasizes user features
- **XGBoost:** Emphasizes temporal patterns
- **CatBoost:** Emphasizes community aggregates
- **Stacking:** Learns when each model's emphasis is reliable

Meta-learner's logistic weights (implicitly learned) allocate trust per class.

9.3 Data Preprocessing Impact on Explainability

Missing value imputation to "Unknown":

- TF-IDF vectorizer creates a feature for unknown patterns
- SHAP shows this feature has near-zero importance
- **Validation:** Imputation strategy doesn't inject spurious signal

10 Results & Performance

10.1 Cross-Validation Results

CROSS-VALIDATION PERFORMANCE (10-Fold Stratified):

Base Models Performance:

- Average Base Model Macro F1 (across 10 folds): **0.9686**
- Standard Deviation: ± 0.0132 (demonstrating reasonable stability)
- All folds consistently achieved F1 scores above 0.95

Ensemble Performance:

- Final Meta-learner Macro F1 (on OOF predictions): **0.9473**

Performance Analysis:

While the meta-learner F1 is slightly lower than the average base model F1, this is not uncommon in stacking scenarios. The meta-learner provides robust, overall performance by learning to trust different base models for different prediction scenarios. The 0.9473 F1 score represents strong generalization on unseen data during cross-validation, demonstrating the ensemble's reliability.

Stability: Low standard deviation indicates robust generalization across all folds.

11 Key Innovations & Why They Matter

11.1 Cyclical Temporal Encoding

Conventional approach: Raw hour, month as linear features → Model struggles with periodicity (11pm vs 1am appear far apart)

Our approach: Sine/cosine encoding → Natural periodicity captured

SHAP validation: Cyclical features appear in top 20, confirming effectiveness.

11.2 Leak-Safe Aggregations

Conventional approach: Global user/community stats → Validation fold sees its own aggregated information → Inflated CV scores

Our approach: Recompute aggregations per fold on training data only → Honest CV estimates

Impact: Prevents false confidence; private LB scores reflect true generalization.

11.3 Stacking with Logistic Regression

Conventional approach: Simple averaging or voting → Loses learned patterns

Our approach: Meta-learner learns base model trustworthiness → Provides robust final predictions by learning trustworthiness patterns of base models across different scenarios

Cost: Minimal—logistic regression is negligible overhead.

11.4 Class Weighting + Early Stopping

Problem: Class imbalance + noisy stopping criteria → Poor minority class performance

Solution:

- Compute per-class weights: $\text{weight} = \frac{N_{\text{total}}}{k \times N_{\text{class}}}$
- Apply via `sample_weight + is_unbalance=True`
- Early stopping on validation set prevents overfitting to majority

Result: Macro F1 (unweighted average) stays high across all classes.

12 Reproducibility & Code Quality

12.1 Deterministic Design

- Fixed random seeds (`SEED=42`) everywhere
- Byte-stable CSV exports
- All preprocessing in code (no manual Excel edits)

12.2 Modular Architecture

1. <code>rebuild_clean_files()</code>	-> Canonical CSVs
2. <code>advanced_features()</code>	-> Feature engineering
3. <code>add_aggregations_cv()</code>	-> Leak-safe aggregation
4. Training loop with 10-fold	-> Base models
5. Stacking meta-learner	-> Final predictions

Each module is self-contained, testable.

12.3 Hyperparameter Justification

Every hyperparameter chosen based on:

- Domain intuition (e.g., 4-hour blocks for diurnal cycles)
- CV stability (e.g., `early_stopping=200` to prevent overfitting)
- Balance (e.g., `num_leaves=127 = 27 - 1`, practical sweet spot)

Not arbitrary tuning.

13 Conclusion

This solution demonstrates that indigenous knowledge, systematically encoded and blended with modern ML, outperforms naive approaches. The 2nd place ranking validates:

- Rigorous data preprocessing prevents leakage
- Thoughtful feature engineering captures domain insights
- Ensemble methods robustly combine diverse perspectives
- SHAP explainability proves the model learns sensible patterns
- Reproducible

For the Zindi competition and beyond: This approach bridges scientific rigor with respect for indigenous ecological expertise, contributing to more inclusive and locally-grounded weather prediction systems.

Submission Category: Advanced Ensemble with Explainability
Author: Med Amin Kharrat